



# Методи чисельної оптимізації

Лекція 2: Автоматичне диференціювання

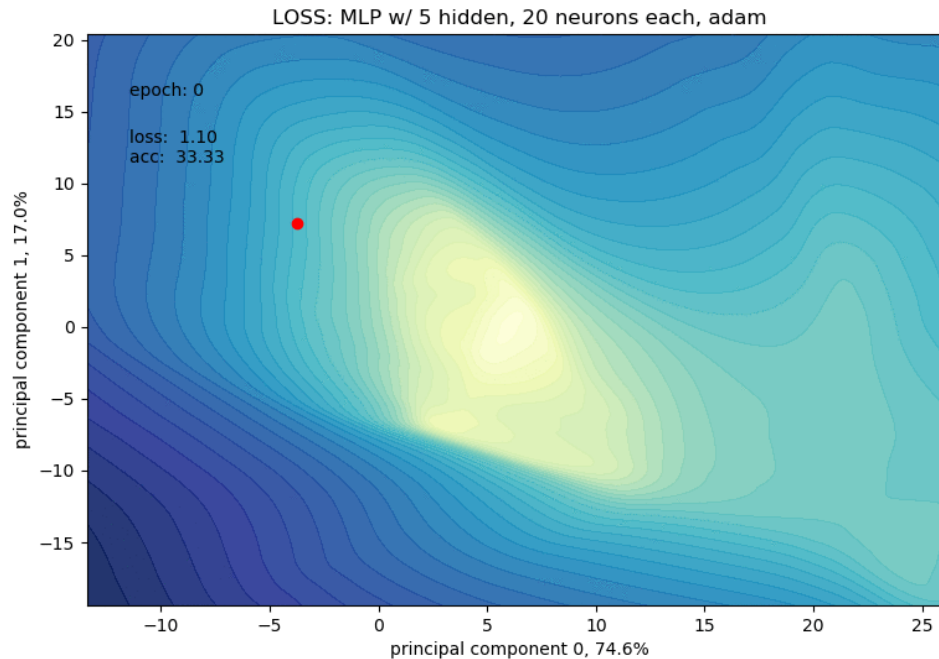
Кочура Юрій Петрович  
[iuriy.kochura@gmail.com](mailto:iuriy.kochura@gmail.com)  
[@y\\_kochura](#)

# Сьогодні

- 🎙 Мотивація
- 🎙 Огляд: похідна, градієнт, яacobіан
- 🎙 Обчислення градієнтів у зворотному поширенні
- 🎙 Імплементация у фреймворках



# Мотивація



## Мотивація

- Алгоритми навчання на основі градієнтів є рушієм глибокого навчання.
- Обчислювати градієнти вручну — справа стомлююча та схильна до помилок. Це швидко стає непрактичним для складних моделей.
- Зміни в моделі вимагають повторного обчислення градієнтів.

## Програми як диференційовані функції

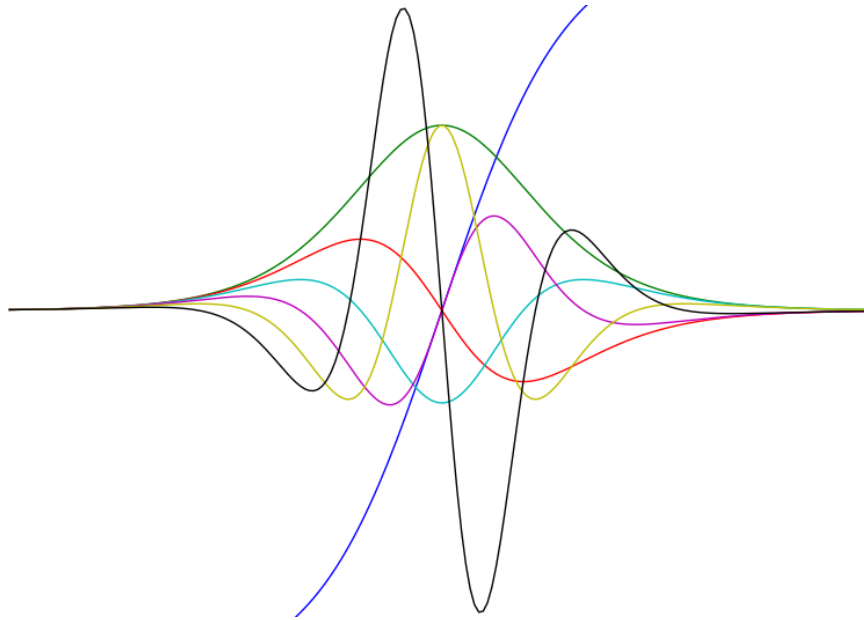
Програма визначається як композиція простих функцій, які ми знаємо, як диференціювати окремо.

```
import jax.numpy as jnp
from jax import grad

def predict(params, inputs):
    for W, b in params:
        outputs = jnp.dot(inputs, W) + b
        inputs = jnp.tanh(outputs)
    return outputs

def loss_fun(params, inputs, targets):
    preds = predict(params, inputs)
    return jnp.mean((preds - targets)**2)

grad_fun = grad(loss_fun)
```



Сучасні фреймворки підтримують похідні вищого порядку.

```
def tanh(x):  
    y = jnp.exp(-2.0 * x)  
    return (1.0 - y) / (1.0 + y)  
  
fp = grad(tanh)  
fpp = grad(grad(tanh))  # друга похідна  
...
```

# Автоматичне диференціювання

Автоматичне диференціювання надає набір методів для оцінки **похідних** функції, заданої **комп'ютерною програмою**.

- $\neq$  символічне диференціювання, яке спрямоване на знаходження деякого зрозумілого людині виразу похідної.
- $\neq$  чисельне диференціювання (наприклад, методом кінцевих різниць) може призвести до помилок округлення.

# Огляд

Похідна, градієнт, яacobian



# Похідна

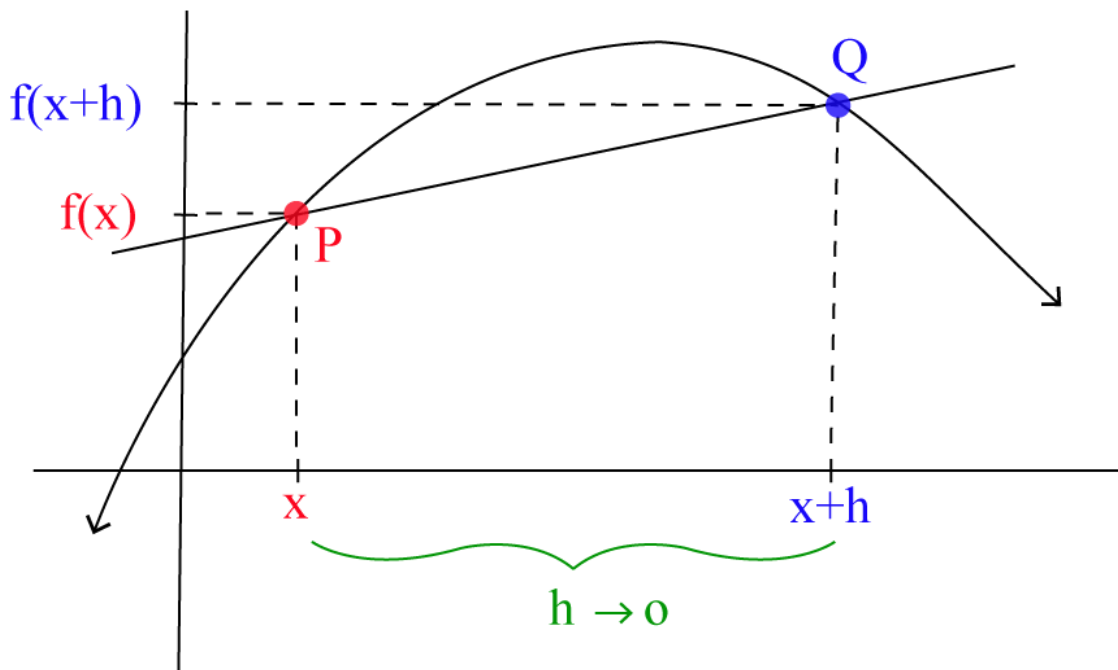
Нехай  $f : \mathbb{R} \rightarrow \mathbb{R}$

Похідна неперервної функції  $f$  визначається як:

$$f'(x) = \frac{\partial f(x)}{\partial x} \triangleq \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h},$$

де

- $f'(x)$  — це позначення похідної за Лагранжом,
- $\frac{\partial f(x)}{\partial x}$  — це позначення похідної за Лейбніцем.



Похідна  $\frac{\partial f(x)}{\partial x}$  показує миттєву швидкість зміни  $f$  у точці  $x$ .

# Градiєнт

Градiєнт неперервної функції  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

$$\nabla f(\mathbf{x}) \triangleq \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n,$$

тобто вектор, який акумулює часткові похідні від  $f$ .

Застосовуючи визначення похідної за координатами, маємо

$$[\nabla f(\mathbf{x})]_j = \frac{\partial f}{\partial x_j}(\mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{e}_j) - f(\mathbf{x})}{h},$$

де  $\mathbf{e}_j$  —  $j$ -й базисний вектор.

# Якобіан

Якобіан  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$

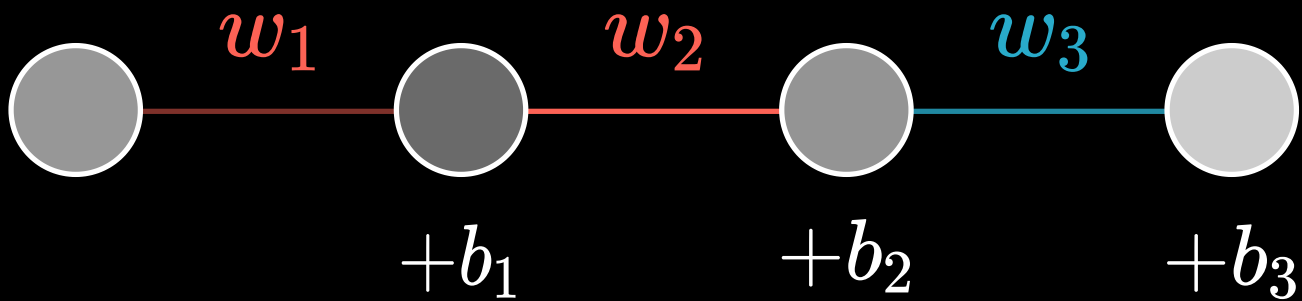
$$\begin{aligned} J_{\mathbf{f}}(\mathbf{x}) &= \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \triangleq \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \dots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}) & \dots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^{m \times n} \\ &= \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1}(\mathbf{x}) & \dots & \frac{\partial \mathbf{f}}{\partial x_n}(\mathbf{x}) \end{bmatrix} \\ &= \begin{bmatrix} \nabla f_1(\mathbf{x})^T \\ \vdots \\ \nabla f_m(\mathbf{x})^T \end{bmatrix} \end{aligned}$$

Для ( $m = 1$ ) якобіан = градієнт.

# Обчислення градієнтів у зворотному поширенні



Почнемо з досить простої мережі, де кожен шар має лише один нейрон.



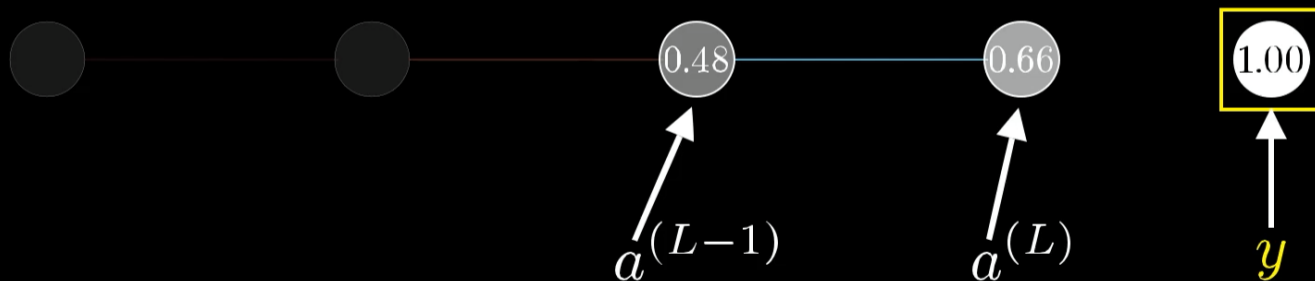




$$C_0(\dots) = (a^{(L)} - y)^2$$

For example:  $(0.66 - 1.00)^2$

Desired  
output



$$a^{(L)} = \sigma(w^{(L)} a^{(L-1)} + b^{(L)})$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

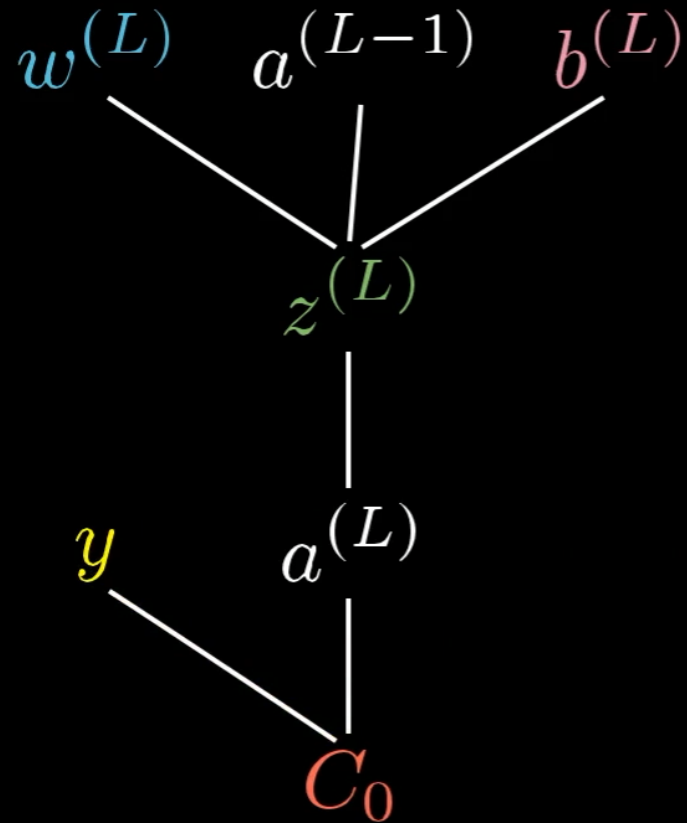
$$a^{(L)} = \sigma(z^{(L)})$$

$$a^{(L)} = \sigma(w^{(L)} a^{(L-1)} + b^{(L)})$$

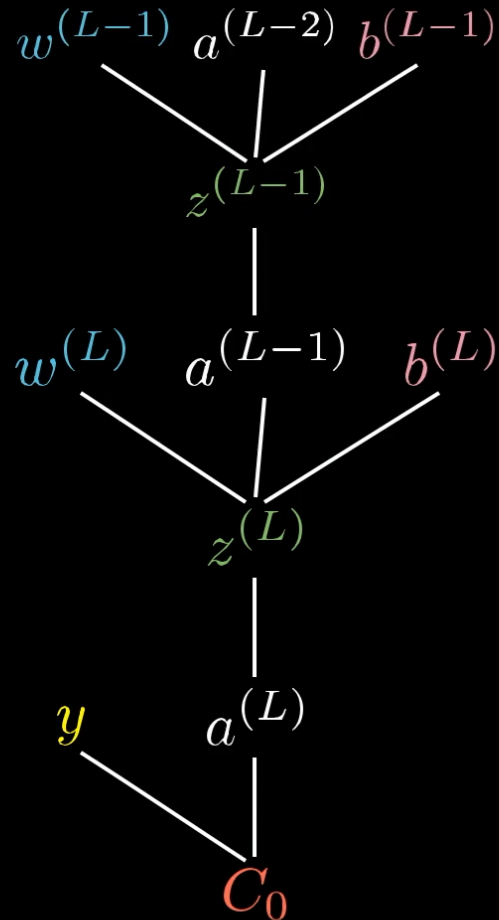
$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

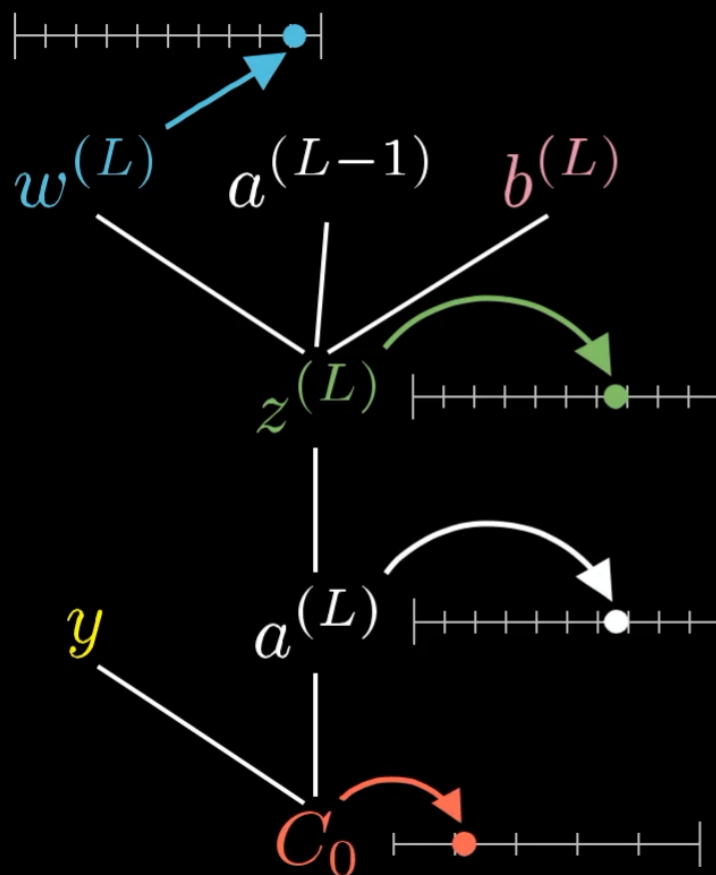
$$a^{(L)} = \sigma(z^{(L)})$$

# Обчислювальний граф



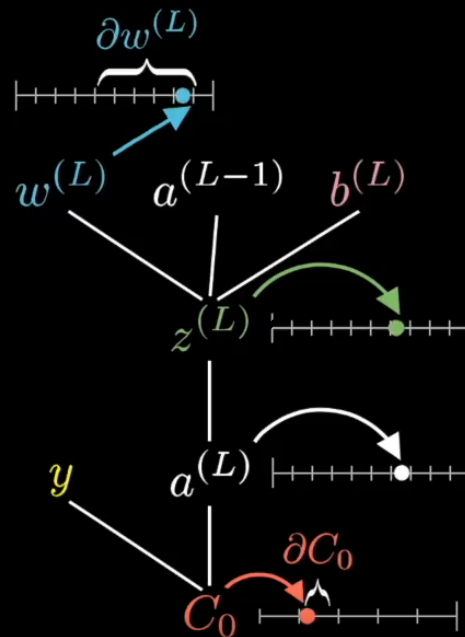
# Обчислювальний граф

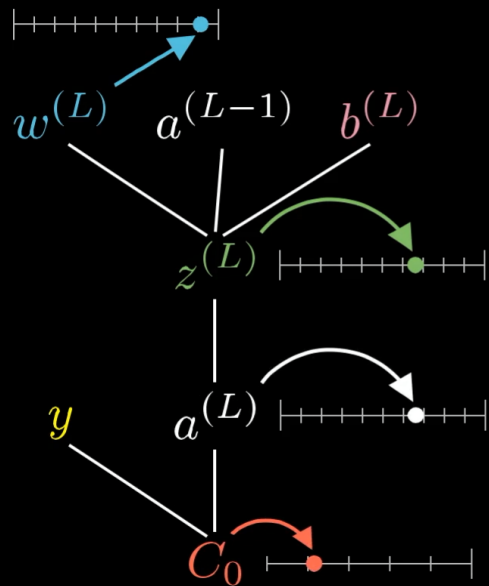




Усі ці змінні просто числа, тому було б добре уявити, що кожна змінна має власну числову шкалу.

Наша перша мета — зрозуміти, наскільки втрата  $C_0$  чутлива до невеликих змін вагових коефіцієнтів  $w^{(L)}$ . Тобто ми хочемо знати похідну  $\frac{\partial C_0}{\partial w^{(L)}}$ .  
Невелике збурення  $w^{(L)}$  має ланцюжок ефектів, які зрештою спричиняють збурення  $C_0$ .





$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial C_0}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial w^{(L)}}$$

Це правило ланцюжка, де множення цих трьох співвідношень дає нам чутливість  $C_0$  до невеликих змін  $w^{(L)}$ .



# Похідні

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)} \quad \longrightarrow \quad \frac{\partial z^{(L)}}{\partial w^{(L)}} = a^{(L-1)}$$

$$a^{(L)} = \sigma(z^{(L)}) \quad \longrightarrow \quad \frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'(z^{(L)})$$

$$C_0 = (a^{(L)} - y)^2 \quad \longrightarrow \quad \frac{\partial C_0}{\partial a^{(L)}} = 2(a^{(L)} - y)$$

$$\boxed{\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial C_0}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial w^{(L)}}}$$

Як бачите, кожна похідна досить проста, якщо знати для якого рівняння шукати похідну.

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)} \quad \longrightarrow \quad \frac{\partial z^{(L)}}{\partial w^{(L)}} = a^{(L-1)}$$

$$a^{(L)} = \sigma(z^{(L)}) \quad \longrightarrow \quad \frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'(z^{(L)})$$

$$C_0 = (a^{(L)} - y)^2 \quad \longrightarrow \quad \frac{\partial C_0}{\partial a^{(L)}} = 2(a^{(L)} - y)$$

$$\boxed{\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial C_0}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial w^{(L)}}}$$

$$\boxed{\frac{\partial C_0}{\partial w^{(L)}} = 2(a^{(L)} - y) \sigma'(z^{(L)}) a^{(L-1)}}$$

$$\frac{\partial C_0}{\partial w^{(L)}} = 2(a^{(L)} - y)\sigma'(z^{(L)})a^{(L-1)}$$

Ця формула показує як вплине значення вагового коефіцієнта  $w^{(L)}$  на втрати мережі. Це лише один дуже конкретний шматок інформації, нам потрібно обчислити ще багато подібних похідних, щоб отримати весь вектор градієнта.

# Цільова функція втрат

$C_0$  — втрати для одного навчального прикладу.

Цільова функція втрат мережі (усереднені втрати для всього набору даних):

$$C = \frac{1}{n} \sum_{k=0}^{n-1} C_k$$

Отже, якщо ми хочемо отримати похідну від  $C$  відносно ваг, нам потрібно взяти середнє значення всіх окремих похідних:

$$\frac{\partial C}{\partial w^{(L)}} = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^{(L)}}$$

Цей вираз визначає, як зміниться цільова втрата мережі, коли ми змінимо ваги останнього шару, тобто  $w^{(L)}$ .

# Обчислення повного градієнта

$$\nabla C \triangleq \begin{bmatrix} \frac{\partial C}{\partial w^{(1)}} \\ \frac{\partial C}{\partial b^{(1)}} \\ \vdots \\ \frac{\partial C}{\partial w^{(L)}} \\ \frac{\partial C}{\partial b^{(L)}} \end{bmatrix}$$

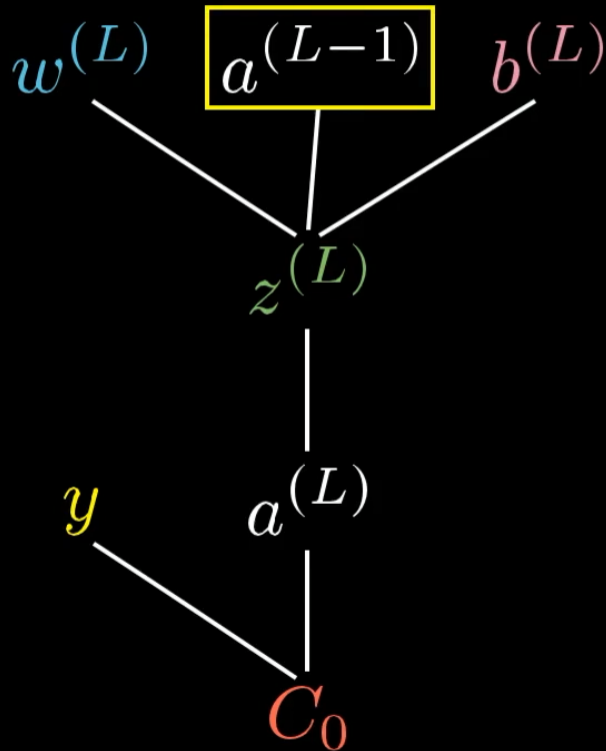
Хороша новина полягає в тому, що чутливість функції витрат до зміни зсуву майже ідентична рівнянню для зміни ваг:

$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial C_0}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial w^{(L)}}$$

$$\frac{\partial C_0}{\partial b^{(L)}} = \frac{\partial C_0}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial b^{(L)}}$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)} \quad \longrightarrow \quad \frac{\partial z^{(L)}}{\partial b^{(L)}} = 1$$

## Вплив вагів та зсуву попередніх шарів



$$\frac{\partial C_0}{\partial a^{(L-1)}} = \frac{\partial C_0}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial a^{(L-1)}}$$

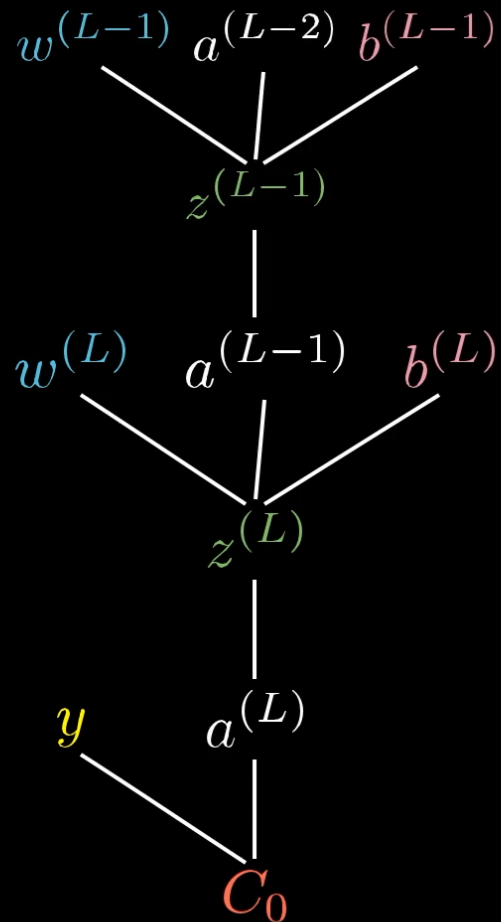
Усі інші ваги та зсуви з попередніх шарів мережі мають менш прямий вплив на цільові втрати.

$$\frac{\partial C_0}{\partial a^{(L-1)}} = \frac{\partial C_0}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial a^{(L-1)}}$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)} \quad \longrightarrow \quad \frac{\partial z^{(L)}}{\partial a^{(L-1)}} = w^{(L)}$$



Активация у передостанньому шарі визначається його власними вагами та зсувом.



$$\frac{\partial C_0}{\partial w^{(L-1)}} = \frac{\partial C_0}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial w^{(L-1)}}$$

$$z^{(L-1)} = w^{(L-1)} a^{(L-2)} + b^{(L-1)} \quad \longrightarrow \quad \frac{\partial z^{(L-1)}}{\partial w^{(L-1)}} = a^{(L-2)}$$

$$a^{(L-1)} = \sigma(z^{(L-1)}) \quad \longrightarrow \quad \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} = \sigma'(z^{(L-1)})$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)} \quad \longrightarrow \quad \frac{\partial z^{(L)}}{\partial w^{(L)}} = a^{(L-1)}, \quad \frac{\partial z^{(L)}}{\partial a^{(L-1)}} = w^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)}) \quad \longrightarrow \quad \frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'(z^{(L)})$$

$$C_0 = (a^{(L)} - y)^2 \quad \longrightarrow \quad \frac{\partial C_0}{\partial a^{(L)}} = 2(a^{(L)} - y)$$

Відстежуючи залежності через наш обчислювальний граф та перемножуючи довгий ряд часткових похідних, тепер ми можемо обчислити похідну функції втрат відносно ваг чи зсуву у будь-якому шарі всієї мережі. Ми просто застосовуємо ту саму ідею правила ланцюжка, якою користувалися раніше!

Оскільки ми можемо отримати будь-яку похідну, ми можемо обчислити весь вектор градієнта. Робота зроблена! Принаймні для цієї мережі.

Відстежуючи залежності через наш обчислювальний граф та перемножуючи довгий ряд часткових похідних, тепер ми можемо обчислити похідну функції втрат відносно ваг чи зсуву у будь-якому шарі всієї мережі. Ми просто застосовуємо ту саму ідею правила ланцюжка, якою користувалися раніше!

Оскільки ми можемо отримати будь-яку похідну, ми можемо обчислити увесь вектор градієнта. Робота зроблена! Принаймні для цієї мережі.

# Імплементація


 PyTorch

 TensorFlow



# Примітиви

Більшість фреймворків, в яких реалізовано автоматичне диференціювання, визначаються набором комбінованих **примітивних** операцій.



Search docs

**TUTORIALS**

- JAX Quickstart
- How to Think in JAX
- The Autodiff Cookbook
- Autobatching log-densities example
- Training a Simple Neural Network, with Tensorflow Datasets Data Loading

**ADVANCED JAX TUTORIALS**

- JAX - The Sharp Bits
- Custom derivative rules for JAX-transformable Python functions
- How JAX primitives work
- Writing custom Jaxpr interpreters in JAX
- Training a Simple Neural Network, with PyTorch Data Loading
- XLA in Python
- Caution: This is a pedagogical notebook covering some low level XLA details, the APIs herein are neither public nor stable!

Docs » Public API: jax package » jax.lax package

Edit on GitHub

## jax.lax package

`jax.lax` is a library of primitives operations that underpins libraries such as `jax.numpy`. Transformation rules, such as JVP and batching rules, are typically defined as transformations on `jax.lax` primitives.

Many of the primitives are thin wrappers around equivalent XLA operations, described by the [XLA operation semantics](#) documentation. In a few cases JAX diverges from XLA, usually to ensure that the set of operations is closed under the operation of JVP and transpose rules.

Where possible, prefer to use libraries such as `jax.numpy` instead of using `jax.lax` directly. The `jax.numpy` API follows NumPy, and is therefore more stable and less likely to change than the `jax.lax` API.

### Operators

<code>abs(x)</code>	Elementwise absolute value: $ x $ .
<code>add(x, y)</code>	Elementwise addition: $x + y$ .
<code>acos(x)</code>	Elementwise arc cosine: $\text{acos}(x)$ .
<code>argmax(operand, axis, index_dtype)</code>	Computes the index of the maximum element.
<code>argmin(operand, axis, index_dtype)</code>	Computes the index of the minimum element.
<code>asin(x)</code>	Elementwise arc sine: $\text{asin}(x)$ .
<code>atan(x)</code>	Elementwise arc tangent: $\text{atan}(x)$ .
<code>atan2(x, y)</code>	Elementwise arc tangent of two variables.
<code>batch_matmul(lhs, rhs[, precision])</code>	Batch matrix multiplication.

## Складання примітивів

Примітивні функції складаються разом у граф, який описує обчислення.

Обчислювальний граф будується:

- **наперед**, з абстрактного синтаксичного дерева програми або за допомогою спеціального API (наприклад, Tensorflow 1), або
- **just in time**, відстежуючи виконання програми (наприклад, Tensorflow Eager, JAX, PyTorch).



# Підсумок

- Автоматичне диференціювання надає набір методів для оцінки похідних функції, заданої комп'ютерною програмою. Це є одним з ключових факторів, що зіграло важливу роль в революції глибокого навчання.
- Автоматичне диференціювання у зворотному режимі ефективніше, коли функція має більше входів, ніж виходів.

# Підсумок

- Автоматичне диференціювання надає набір методів для оцінки **похідних** функції, заданої **комп'ютерною програмою**. Це є одним з ключових факторів, що зіграло важливу роль в революції глибокого навчання.
- Автоматичне диференціювання у зворотному режимі ефективніше, коли функція має більше входів, ніж виходів.



# Література

- Mathieu Blondel, [Automatic differentiation](#), 2020.
- Gabriel Peyré, [Course notes on Optimization for Machine Learning](#), 2020.